
INTERNATIONAL CONFERENCE ON MODELLING, OPTIMISATION AND
COMPUTING (ICMOC2012)

Lung Nodule Segmentation through Unsupervised Clustering Models

S.Sivakumar^a, C.Chandrasekar^b, a *

^aResearch Scholar, Department of Computer Science, Periyar University, Salem, Tamilnadu, India-636 011.

^bReader, Department of Computer Science, Periyar University, Salem, Tamilnadu, India-636 011.

Abstract

Image processing is an essential technique for analyzing images. The important part of image processing is image segmentation. Segmentation is a task of grouping pixels based on similarity. In medical image analysis, segmentation is very important phase. In this paper Possibilistic Clustering models, Fuzzy Clustering models and a new approach called Possibilistic-Fuzzy based clustering model are discussed. Experiments are carried out on bench mark medical images to examine the performance of the above techniques. The results are compared with various validation measures to explore the accuracy of our proposed approach.

© 2012 Published by Elsevier Ltd. Selection and/or peer-review under responsibility of Noorul Islam Centre for Higher Education. Open access under [CC BY-NC-ND license](#).

Keywords: Possibilistic clustering, Fuzzy clustering, Possibilistic-Fuzzy clustering, Lungs nodule segmentation, Validity measures.

1. Introduction

1.1 Lung Cancer

Lung cancer is the primary cause of tumor deaths for both sexes in most countries. There are four stages of lung cancer from I to IV with rising gravity. If the cancer is detected at stage I and it has no more 30 mm in diameter, then there is about 67% survival rate, and only less than 1% chance left for stage IV. Thus it is concluded that early detection and treatment at stage 1 have high survival rate. But unfortunately, lung cancer is usually detected late due to the lack of symptoms in its early stages. This is the reason why lung screening programs have been investigated to detect pulmonary nodules: they are

Corresponding author. Tel.: +91-9865103530.

E-mail address: ssivakkumarr@yahoo.com.

small lesions which can be calcified or not, almost spherical in shape or with irregular borders. The nodule definition for thoracic CT of the Fleischner's Society is "a round opacity, at least moderately well margined and no greater than 3 cm in maximum diameter" [7]. Approximately 40% of lung nodules are malignant, that is, are cancerous; the rest is usually associated with infections. Because malignancy depends on many factors, such as patient age, nodule shape, doubling time, presence of calcification [8], after the initial nodule detection further exams are necessary to obtain a diagnosis. In computer vision, segmentation refers to the process of partitioning a digital image into multiple regions or sets of pixels. Each of the pixels in a region is similar with respect to some characteristic or computed property, such as color, intensity, or texture. Adjacent regions are significantly different with respect to the same characteristics [9][10][11]. Early diagnosis has an important prognostic value and has a huge impact on treatment planning [1]. As nodules are the most common sign of lung cancer, nodule detection in CT scan images is a main diagnostic problem. Conventional projection radiography is a simple, cheap, and widely used clinical test. Unfortunately, its capability to detect lung cancer in its early stages is limited by several factors, both technical and observer-dependent. Lesions are relatively small and usually contrast poorly with respect to anatomical structure. This partially explains why radiologists are commonly credited with low sensitivity in nodule detection, ranging from 60 to 70%. A thorough review of the drawbacks affecting conventional chest radiography is given, for example, by Woodring [2]. However, several long-term studies carried out in the 1980s using large clinical data sets have shown that up to 90% of nodules may be correctly perceived retrospectively [3], [4]. In addition, detection sensitivity can be increased to more than 80% in the case of a double radiograph reading by two radiologists. Furthermore, sensitivity is expected to increase with the widespread use of digital radiography systems which are characterized by an extended dynamic range and have a better contrast resolution than conventional film radiography. In view of this, the availability of efficient and effective computer-aided diagnosis (CAD) systems is highly desirable [5], as such systems are usually conceived to provide the physician with a second opinion [6] so as to focus his/her attention on suspicious image zones, playing the role of a "second reader."

1.2 Image Segmentation

Image segmentation remains one of the major challenges in image analysis, since image analysis tasks are often constrained by how well previous segmentation is accomplished. In particular, many existing image segmentation algorithms fail to provide satisfactory results when the boundaries of the desired objects are not clearly defined by the image intensity information. Image segmentation has a large number of applications in several fields. Its importance concerns basically in distinguishing different objects or regions inside the image. The process of automating image segmentation is, however, very complicated [12]. This is often unfeasible, not only because of technical limitations but because there is no indication of what is meaningful for an intended user. Natural scenes, for example, are rich in details and tonal variations, which are better detected when the algorithm considers additional features. Therefore, texture based methods tend to be more efficient than tone based methods. One of the most used approaches to image segmentation is based on region detection. Region based methods consist in localizing areas by analyzing pixels statistics. Several authors have worked in finding descriptors and features for texture identification. Among these, Haralick features [13] are the most widely used. In his work, Haralick [13] suggested the use of gray-level co-occurrence matrices (GLCM) to extract texture features from an image. Since then, GLCMs became widely used for image texture features extraction in many types of applications.

In the field of medical diagnosis an extensive diversity of imaging techniques is presently available, such as radiography, computed tomography (CT) and magnetic resonance imaging (MRI). In recent times,

Computed Tomography (CT) is the most effectively used for diagnostic imaging examination for chest diseases such as lung cancer, tuberculosis, pneumonia and pulmonary emphysema. The volume and the size of the medical images are progressively increasing day by day. Therefore it becomes necessary to use computers in facilitating the processing and analyzing of those medical images. Even though the original FCM algorithm yields good results for segmenting noise free images, it fails to segment images corrupted by noise, outliers and other imaging artifact. Medical image segmentation is an indispensable step for most successive image analysis tasks. This paper presents various image segmentation approaches based on unsupervised clustering models.

2 Clustering Models

Cluster analysis is a method for clustering a data set into groups of similar individuals. It is an approach towards unsupervised learning as well as one of the major techniques in pattern recognition. The conventional (hard) clustering methods restrict each point of the data set to exactly one cluster. Since Zadeh [14] proposed fuzzy sets that produced the idea of partial membership described by a membership function, fuzzy clustering has been widely studied and applied in a variety of key areas [15][16][17]. In the literature on fuzzy clustering, the fuzzy c-means (FCM) clustering algorithm, proposed by Dunn [18] and extended by Bezdek [15], is the most well-known and used method.

Although FCM is a very useful clustering method, its memberships do not always correspond well to the degrees of belonging of the data, and it may be inaccurate in a noisy environment [19]. To improve this weakness of FCM, and to produce memberships that have a good explanation of the degrees of belonging for the data, Krishnapuram and Keller [19] created a possibilistic approach to clustering which used a possibilistic type of membership function to describe the degree of belonging. They showed that algorithms with possibilistic memberships are more robust to noise and outliers than FCM. The possibilistic clustering approach has also been applied in shell clustering, boundary detection, surface and function approximations [20] [21] [22].

It is necessary to pre-assume the number c of clusters for these hard, fuzzy and possibilistic clustering algorithms. In general, the cluster number c should be unknown. The problem for finding an optimal c is usually called cluster validity. Once the partition is obtained by a clustering method, the validity function can help us to validate whether it accurately presents the structure of the data set or not.

2.1 The FCM clustering algorithm

In the unsupervised learning literature, the FCM is the best-known fuzzy clustering method. The FCM is an iterative algorithm using the necessary conditions for a minimizer of the FCM objective function J_{FCM} with

$$J_{FCM}(\mu, a) = \sum_{i=1}^c \sum_{j=1}^n \mu_{ij}^m \|x_j - a_i\|^2, m > 1 \quad (1)$$

Where $\mu = \{\mu_1, \dots, \mu_c\}$ with the membership function μ_i defined as $\mu_{ij} = \mu_i(x_j)$ is a fuzzy c -partition and $a = \{a_1, \dots, a_c\}$ is the set of c cluster centers. The necessary conditions for a minimizer (μ, a) of J_{FCM} are the following update equations:

$$\mu_{ij} = \frac{\|x_j - a_i\|^{\frac{2}{m-1}}}{\sum_{k=1}^c \frac{\|x_j - a_k\|^{\frac{2}{m-1}}}{\|x_j - a_i\|^{\frac{2}{m-1}}}}, \quad i = 1, \dots, c, \quad j = 1, \dots, n \quad (2)$$

and

$$a_i = \frac{\sum_{j=1}^n \mu_{ij}^m x_j}{\sum_{j=1}^n \mu_{ij}^m}, \quad i = 1, \dots, c. \quad (3)$$

The weighting exponent m is called the fuzzifier which can have an influence on the clustering performance of FCM [23].

Fuzzy c-means clustering algorithm:

Initialize $a_i^{(0)}$, $i=1, \dots, c$ and set $\varepsilon > 0$; set iteration counter $l=0$;

Step1: Compute $\mu_{ij}^{(l+1)}$ using eq. (2).

Step2: Compute $a_i^{(l+1)}$ using eq. (3).

Increment l ; until $\max_i \|a_i^{(l+1)} - a_i^{(l)}\| < \varepsilon$

2.2 The Possibilistic Clustering algorithm

Although FCM is a very useful clustering method, its memberships do not always correspond well to the degree of belonging of the data, and may be inaccurate in a noisy environment [19]. To improve this weakness of FCM, and to produce memberships that have a good explanation for the degree of belonging for the data, Krishnapuram and Keller [19] relaxed the constrained condition $\sum_{i=1}^c \mu_i(x) = 1$ of the fuzzy c -partition $\{\mu_1, \dots, \mu_c\}_F$ in FCM to obtain a possibilistic type of membership function with

$$\mu_1, \dots, \mu_c : P = \{ \mu_1, \dots, \mu_c \mid \max_i \mu_i(x) > 0, \quad i \} \quad (4)$$

We may call the memberships μ_1, \dots, μ_c in Eq. (4) as the possibilistic c -memberships. To avoid the trivial solutions, Krishnapuram and Keller [19] added a constraining term to FCM and proposed the following possibilistic clustering objective function:

$$J(\mu, A) = \sum_{i=1}^c \sum_{j=1}^n \mu_{ij}^m \|x_j - a_i\|^2 + \sum_{i=1}^c \sum_{j=1}^n \eta_i (1 - \mu_{ij})^m \quad (5)$$

$$\eta_i = K \frac{\sum_{j=1}^n \mu_{ij}^m \|x_j - a_i\|^2}{\sum_{j=1}^n \mu_{ij}^m} \quad (6)$$

$$\mu_{ij} = \frac{1}{1 + \left[\frac{\|x_j - a_i\|^2}{\eta_i} \right]^{\frac{1}{m-1}}}}, \quad 1 \leq i \leq c; \quad (7)$$

Where $K \in (0, \infty)$ was typically chosen to be one. They then created a possibilistic approach to clustering which used a possibilistic c -membership of Eq. (4) to describe the degree of belonging on the basis of the objective function (5).

2.3 Fuzzy-Possibilistic Clustering Model

FPCM algorithm was proposed by N.R.Pal, K.Pal, and J.C.Bezdek [25] and it includes both possibility and membership values. FPCM model can be seen as below:

$$\min\{J_{M,\eta}(U,T,V;X)\} = \sum_{i=1}^c \sum_{k=1}^n (u_{ik}^m + t_{ik}^\eta) \|x_j - c_i\|^2 \quad (8)$$

Where U is membership matrix, T is possibilistic matrix, and V is the resultant cluster centers, c and n are cluster number and data point number respectively.

$$\mu_{ij} = \frac{1}{\sum_{j=1}^n \left(\frac{\|x_j - c_i\|^2}{\|x_k - c_j\|^2} \right)^{\frac{2}{m-1}}}, 1 \leq i \leq c; 1 \leq k \leq n \quad (9)$$

$$t_{ik} = \frac{1}{\sum_{j=1}^n \left(\frac{\|x_j - c_i\|^2}{\|x_k - c_j\|^2} \right)^{\frac{2}{\eta-1}}}, 1 \leq i \leq c; 1 \leq k \leq n \quad (10)$$

$$v_i = \frac{\sum_{k=1}^n (u_{ik}^m + t_{ik}^\eta) x_k}{\sum_{k=1}^n (u_{ik}^m + t_{ik}^\eta)}, 1 \leq i \leq c \quad (11)$$

The above equations show that membership u_{ik} is affected by all c cluster centers, while possibility t_{ik} is affected only by the i -th cluster center c . The possibilistic term distributes the t_{ik} with respect to all n data points, but not with respect to all c clusters. So, membership can be called relative typicality, it measures the degree to which a point belongs to one cluster relative to other clusters and is used to crisply label a data point. And possibility can be viewed as absolute typicality, it measures the degree to which a point belongs to one cluster relative to all other data points, it can reduce the effect of outliers. Combining both membership and possibility can lead to better clustering result.

2.4 Proposed Weighted Fuzzy-Possibilistic Clustering Model

The objective function of the proposed Weighted Fuzzy-Possibilistic Clustering can be formulated as follows:

$$J_{MFPCM} = \sum_{i=1}^c \sum_{j=1}^n \left[\mu_{ij}^{2m} w_{ji}^m \|x_j - c_i\|^{2m} + t_{ij}^{2\eta} w_{ji}^m \|x_j - c_i\|^{2m} \right] \quad (12)$$

Where

$$\mu_{ij} = \frac{1}{\sum_{k=1}^n \left[\frac{\|x_j - c_i\|^{2m}}{\|x_k - c_j\|^{2m}} \right]^{\frac{2m}{m-1}}} \quad (13)$$

$$t_{ij} = \frac{1}{\sum_{k=1}^C \frac{\|x_j - c_k\|^{2m}}{w_{kj}^{2m}}} \quad (14)$$

$$v_i = \frac{\sum_{j=1}^N \left(\frac{\|x_j - c_i\|^{2m}}{w_{ij}^{2m}} \right) x_j}{\sum_{j=1}^N \left(\frac{\|x_j - c_i\|^{2m}}{w_{ij}^{2m}} \right)} \quad (15)$$

$$w_k = \prod_{y=1}^m \exp(-h \times \|x_k - x_y\| / \sigma) \quad (16)$$

Where h is a resolution parameter and σ = standard deviation of input data. The proposed Weighted Fuzzy-Possibilistic clustering yields better results when compare with the existing FPCM, PCM, FCM models.

2.5 Validation Measures:

Cluster validity refers to the problem whether a given fuzzy partition fits to the data all. The clustering algorithm always tries to find the best fit for a fixed number of clusters and the parameterized cluster shapes. However this does not mean that even the best fit is meaningful at all. Either the number of clusters might be wrong or the cluster shapes might not correspond to the groups in the data, if the data can be grouped in a meaningful way at all. Different scalar validity measures have been proposed in the literature, none of them is perfect by oneself. In order to evaluate the segmentation validity we use the Partition Coefficient (PC), Classification Entropy (CE), Partition Index (SC), Separation Index(S), and Xie-Beni Index (XB).

2.5.1 Partition Coefficient (PC): measures the amount of “overlapping” between clusters. It is defined by Bezdek[15] as follows:

$$PC(C) = \frac{1}{N} \sum_{i=1}^C \sum_{j=1}^N (\mu_{ij})^2 \quad (17)$$

Where μ_{ij} is the membership of data point j in cluster i . the disadvantage of PC is lack of direct connection to some property of the data themselves. The optimal number of clusters is at the maximum value.

2.5.2 Classification Entropy (CE): it measures the fuzziness of the cluster partition only, which is similar to the Partition Coefficient.

$$CE(C) = -\frac{1}{N} \sum_{i=1}^C \sum_{j=1}^N \mu_{ij} \log(\mu_{ij}) \quad (18)$$

2.5.3 Partition Index (SC): it is the ratio of sum of compactness and separation of clusters. It is a sum of individual cluster validity measures normalized through division by the fuzzy cardinality of each cluster [12].

$$SC(C) = \frac{\sum_{i=1}^C \sum_{j=1}^N (\mu_{ij})^m \|x_j - v_i\|^2}{N \sum_{k=1}^C \sum_{j=1}^N \mu_{kj}^m \|x_j - v_k\|^2} \quad (19)$$

SC is useful when comparing different partitions having equal number of clusters. A lower value of SC indicates a better partition.

2.5.4 Separation Index(S): on the contrary of Partition Index (SC), the separation index uses a minimum-distance separation for partition validity [12].

$$S(C) = \frac{\sum_{j=1}^N (\mu_{ij})^m \|x_j - v_i\|^2}{N m \min_{i,k} \|x_j - v_i\|^2} \quad (20)$$

2.5.5 Xie-Beni Index (XB): it aims to quantify the ratio of the total variation within clusters and the separation of clusters [24].

$$XB(C) = \frac{\sum_{j=1}^N (\mu_{ij})^m \|x_j - v_i\|^2}{N m \min_{i,k} \|x_j - v_i\|^2} \quad (21)$$

The optimal number of clusters should minimize the value of the index.

3 Materials and methods

3.1 LIDC Dataset

The Lung Image Database Consortium image collection (LIDC-IDRI) consists of diagnostic and lung cancer screening thoracic CT scans with marked-up annotated lesions. It is a web-accessible international resource for development, training, and evaluation of computer-assisted diagnostic (CAD) methods for lung cancer detection and diagnosis. The LIDC-IDRI collection contained on The Cancer Imaging Archive (TCIA) is the complete data set of all 1,010 patients which includes all 399 pilot CT cases plus the additional 611 patient CTs and all 290 corresponding chest x-rays. The lungs image data, nodule size list and annotated XML file documentations can be downloaded from the National Cancer Institute website: <https://wiki.cancerimagingarchive.net/display/Public/LIDC-IDRI>.

3.2 Preprocessing

All the lungs CT scan images are in the format of DICOM with the size of 512X512. To enhance the CT scan images, 3X3 window based median filter was applied to remove the noise.

3.3 Feature Extraction

After enhancing the CT scan image 4X4 windowed co-occurrence matrices were formed with distance=1 and angle=0°. From those co-occurrence matrices 14 Haralick features were calculated. These features are taken as input data for the FCM, PCM, FPCM and WFPCM algorithms. The clustered data validated against the following indices: Partition Coefficient (PC), Classification Entropy (CE), Partition Index (SC), Separation Index(S), and Xie-Beni Index (XB).

4 Results and Discussions

For the experiment we taken 10 CT scan lung cancer affected images of size 512X512 with number of cluster is 5, with fuzziness value m=2 and e=0.0001. Table 1 shows the various validation measures are compared with the FCM, PCM, FPCM and WFPCM.

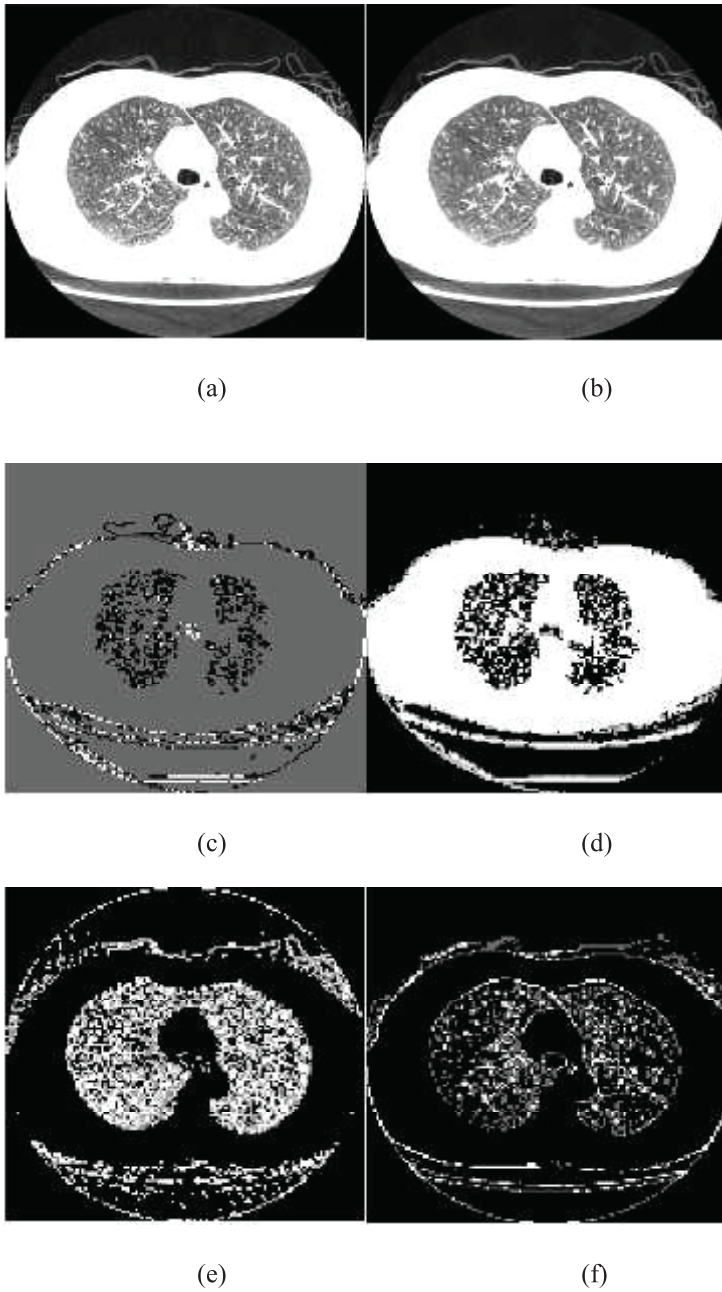


Fig. 1. (a) Original CT scan lungs image (b) preprocessed image (c) FCM applied image (d) PCM applied image (e) FPCM applied image (f) WFPCM applied image

Table1: Cluster validity indices for 30 different CT scan Lungs images with different cluster models (mean value of the indices)

Clustering method	Partition Co-efficient(PC)	Classification Entropy (CE)	Partition Index (SC)	Separation Index (S)	Xie-Beni Index(XB)	No.of. Iteration
FCM	4.0529e-009	2.2781e-008	1.2374e-005	1.3874e-009	24.63	135
PCM	4.3257e-009	2.5381e-008	1.2046e-005	2.2471e-009	24.81	127
FPCM	4.8752e-009	4.3707e-008	1.1308e-005	4.9624e-009	19.42	83
WFPCM	5.2307E-009	6.8506E-009	1.0435E-004	9.7411E-009	14.68	49

From Table1, we can see that the average time spent in WFPCM is far less than that of FCM, PCM and FPCM.

5. Conclusion

In this paper we propose a new algorithm called weighted Fuzzy-Possibilistic C-Means (WFPCM), which is based on adding weight component to both the membership and possibilistic value. A comparison of FCM, PCM, FPCM and WFPCM shows that clustering of normal FCM will be less sensitive to outliers. Finally, the experimental results with the LIDC Datasets shows that WFPCM deals with the amount of noise data, and produces less clustering time and better clustering accuracy.

References

- [1] C. Society, Cancer Facts and Figures 2001. Atlanta, GA: American Cancer Society, 2001.
- [2] J. Woodring.: "Pitfalls in the radiologic diagnosis of lung cancer," AJR, 1990, p.1165–1175.
- [3] J. Muhm, W. Miller, R. Fontana, D. Sanderson, and M. Uhlenhopp.: "Lung cancer detected during a screening program using four-month chest radiographs," Radiology, 1983, vol. 148, p.609–615,.
- [4] N. Hayabuchi, W. Russel, J. Murakami, and H. Nishitani.: "Screening for lung cancer in a fixed population by biennial chest radiography," Radiology, 1983, vol. 148, p.369–373,.
- [5] Van Ginneken, B. M. terHaarRomeny, and M. Viergever.: "Computer-aided diagnosis in chest radiography: A survey," IEEE Trans. Med. Imag., 2001, vol. 20, p.1228–1241.
- [6] T. Kobayashi, X.-W. Xu, H. MacMahon, C. Metz, and K. Doi.: "Effect of a computer-aided diagnosis scheme on radiologists's performance in detection of lung nodules on radiographs," Radiology, 1996, vol. 199, p.843–848.
- [7] J.H.Austin,N.L.Mueller,P.J.Friedman,etal., "Glossary of terms for CT of the lungs: recommendation of theNomenclature Committee of the FleischnerSociety", Radiology1996, vol.200,p.327-331
- [8] <http://www.nlhhep.org>
- [9] AristofanesC. Silva, Paulo Cezar, Marcello Gattas, "Diagnosis of Lung Nodule using GiniCoefficient and skeletonizationin computerized Tomography images",ACMSymposiumonAppliedComputingMarch2004.
- [10] AymanEl-Baz, AlyA. Farag, Robert Falk, Renato La Rocca, "Detection, VisualizationandidentificationofLung Abnormalities in Chest Spiral CT Scan:Phase-I", International Conference on Biomedical Engineering,Cairo, Egypt,12-01-2002
- [11] N.A.Memon,A.M.Mirza,S.A.M.Gilani, "Segmentation ofLungsfromCTScanImgesforEarlyDiagnosisofLungCancer", Proceedings of World Academy of Science, Engineering andTechnology,aug2006 ,vol14.
- [12] Robert M. Haralick, K Shanmugam and Its'Hak Dinstein, "Textural Features for Image Classification." IEEE Transactions on Systems, Man, and Cybernetics, 1979.
- [13] V. Mezaris, I. Kompatsiaris, and M. G. Strintzis, "Still image segmentation tools for content-based multimedia applications," International Journal of pattern recognition and artificial intelligence, Jun. 2004, vol. 18, no. 4, p.701–725.
- [14] L.A. Zadeh, Fuzzy sets, Inf. Control 8, 1965, p.338–353.
- [15] J.C. Bezdek, Pattern Recognition with Fuzzy Objective Function Algorithms, Plenum Press, New York, 1981.
- [16] M.S. Yang, A survey of fuzzy clustering, Math. Comput. Model, 1993,vol. 18, p.1–16.
- [17] F. Höppner, F. Klawonn, R. Kruse, T. Runkler, Fuzzy Cluster Analysis: Methods for Classification Data Analysis and Image Recognition, Wiley, New York, 1999.
- [18] J.C. Dunn, A fuzzy relative of the ISODATA process andits use in detecting compact, well-separated clusters, Journal of Cybernetics, 1974, vol. 3, p.32–57.
- [19] R. Krishnapuram, J.M. Keller, A possibilistic approach to clustering, IEEE Trans. Fuzzy Syst., 1993, vol.1 p. 98–110.

- [20] R. Krishnapuram, H. Frigui, O. Nasraoui, Fuzzy and possibilistic shell clustering algorithm and their application to boundary detection and surface approximation, *IEEE Trans. Fuzzy Syst.*, 1995, vol.3 p. 29–60.
- [21] T.A. Runkler, J.C. Bezdek, Function approximation with polynomial membership functions and alternating cluster estimation, *Fuzzy Sets Syst.*, 1999, vol. 101 p. 207–218.
- [22] T.A. Runkler, J.C. Bezdek, Alternating cluster estimation: a new tool for clustering and function approximation, *IEEE Trans. Fuzzy Syst.*, 1999, vol.7 p.377–393.
- [23] N.R. Pal, J.C. Bezdek, On cluster validity for fuzzy c-means model, *IEEE Trans. Fuzzy Syst.*, 1995, vol.1 p.370–379.
- [24] Weiling Cai, Songcan Chen, Daoqiang Zhang, “Fast and Robust Fuzzy C-Means clustering algorithms incorporating local information for image segmentation”, *Pattern Recognition*, 2007.
- [25] N.R. Pal, and J.C. Bezdek, “A mixed c-means clustering model”, In *IEEE Int. Conf. Fuzzy Systems*, Spain, 1997, p.11-21.